

Revised version

Original article: *Physica A* 444 (2016) 928–939, DOI: <http://dx.doi.org/10.1016/j.physa.2015.10.048>

Corrigendum: *Physica A* 447 (2016) 569–570, DOI: <http://dx.doi.org/10.1016/j.physa.2015.12.010>

A pathway-based network analysis of hypertension-related genes

Huan Wang^{a,b}, Jing-Bo Hu^{b,c}, Chuan-Yun Xu^b, De-Hai Zhang^d, Qian Yan^e,
Ming Xu^{b,f}, Ke-Fei Cao^{b,*}, Xu-Sheng Zhang^{g,h}

^a*School of Computer Science and Technology, Baoji University of Arts and Sciences, Baoji, Shaanxi 721016, China*

^b*Center for Nonlinear Complex Systems, Department of Physics, School of Physics Science and Technology, Yunnan University, Kunming, Yunnan 650091, China*

^c*School of Electronic and Electrical Engineering, Baoji University of Arts and Sciences, Baoji, Shaanxi 721016, China*

^d*School of Software, Yunnan University, Kunming, Yunnan 650091, China*

^e*School of Physics and Electronic Science, Chuxiong Normal University, Chuxiong, Yunnan 675000, China*

^f*School of Mathematical Sciences, Kaili University, Kaili, Guizhou 556011, China*

^g*Modelling and Economics Unit, Centre for Infectious Disease Surveillance and Control, Public Health England, 61 Colindale Avenue, London NW9 5EQ, UK*

^h*Medical Research Council Centre for Outbreak Analysis and Modelling, Department of Infectious Disease Epidemiology, School of Public Health, Imperial College London, Norfolk Place, London W2 1PG, UK*

Abstract

Complex network approach has become an effective way to describe interrelationships among large amounts of biological data, which is especially useful in finding core functions and global behavior of biological systems. Hypertension is a complex disease caused by many reasons including genetic, physiological, psychological and even social factors. In this paper, based on the information of biological pathways, we construct a network model of hypertension-related genes of the salt-sensitive rat to explore the interrelationship between genes. Statistical and topological characteristics show that the network has the small-world but not scale-free property, and exhibits a modular structure, revealing compact and complex connections among these genes. By the threshold of integrated centrality larger than 0.71, seven key hub genes are found: *Jun*, *Rps6kb1*, *Cyts*, *Creb3l2*, *Cdk4*, *Actg1* and *RT1-Da*. These genes should play an important role in hypertension, suggesting that the treatment of hypertension should focus on the combination of drugs on multiple genes.

Keywords: Complex network; Hypertension; Hub gene; Pathway; Modular structure

1. Introduction

The study of complex systems clearly shows that the global behavior of systems is determined by their structure rather than by the properties of their individual parts. The complex network approach has become a powerful tool for studying complex systems, and the global properties of systems are usually studied by abstracting individual elements of systems into nodes and reducing interactions between elements to edges between nodes [1, 2, 3, 4, 5]. Such an approach has been widely applied to understanding gene functions in biological and medical research [6, 7, 8, 9, 10].

*Corresponding author. Tel.: +86 871 65031605.

Email addresses: hwang227@126.com (Huan Wang), kfcao163@163.com (Ke-Fei Cao), xu-sheng.zhang@phe.gov.uk (Xu-Sheng Zhang)

Essential hypertension, which accounts for about 90%–95% of all cases of hypertension [11], is a disease caused by long-term interaction between genetic and environmental factors, and salt is one of the important environmental factors [12]. The blood pressure response to salt loading or salt restriction is heterogeneous among individuals, which is known as salt sensitivity [13, 14, 15]. Salt sensitivity is the genetic susceptibility of individual blood pressure response to salt, and is an intermediate phenotype of essential hypertension [16, 17]. The people who suffer the salt-sensitive (SS) hypertension account for about 50% of hypertensive patients [15]. Although the clinical research and treatment of hypertension have improved dramatically [18, 19, 20], its molecular mechanisms and pathologies involved are still difficult to ascertain.

Many omic data have been obtained and become available through advanced high-throughput technologies, which provide the basis for studying the relationship of biological data by network approach [7, 10]. Various biomolecular networks have been constructed to discover essential functions and mechanisms of biological phenomena [6, 8, 9]. For instance, Censi et al. studied the gene regulatory networks induced in heart tissue by atrial fibrillation [21]. Demicheli and Coradini analyzed breast cancer behavior using gene regulatory networks [22]. Therefore, it is of significance to understand hypertension disease at system level using the complex network approach.

In our previous study [23], we constructed a hypertension-related gene co-expression network by focusing on the analysis of gene expression data (GED) [24] among the Dahl SS rat [25, 26] and two consomic rat strains [27, 28], where the 335 nodes are individual genes and the connections are derived from the expression correlations. This is a theoretical analysis based on GED to determine the key hub genes (nodes) and explore the relationship between these hub genes and hypertension. However, to get more biologically relevant information about hypertension, a pathway-based gene network should also be constructed using the actual biological correlations.

In the present work, we attempt to study the genes that are involved in SS hypertension based on the information of biological pathways. A biological pathway is a series of actions among molecules in a cell that leads to a certain product or a change in a cell [29, 30]. Such a pathway can trigger the assembly of new molecules or turn genes on and off. Since biological pathways such as metabolic and signal transduction pathways can directly be viewed as interconnected processes of molecular species in the cell [29, 30]; therefore, constructing a pathway-based gene network could help to disentangle the actual biological interactions between genes.

In this paper, we will construct the network model of hypertension-related genes according to whether these genes are involved in the same pathways in the KEGG¹ database. Network approach will be employed to investigate the possible relations between network structure and hypertension-related genes based on these data. Through calculating several statistical indices and analyzing topological characteristics of the network, we find that the pathway-based gene network exhibits the small-world but not scale-free property. Meanwhile, the network also exhibits a modular structure: the nodes of the network can be properly divided into groups within which the nodes are highly connected, but between which they are much less connected. The modular structure analysis can visualize the weak connections of the network, and thus help us to study drug targets of hypertension. The results from this paper and the analysis in Ref. [23] would complement each other.

The rest of this paper is organized as follows. In Section 2, we introduce the data source and construct the pathway-based gene network model of hypertension. In Section 3, we analyze the statistical and topological characteristics of the gene network. The modular structure of the network is presented in Section 4, while Section 5 presents summary and concluding remarks.

2. Data source and network construction

The Dahl SS rat, proposed by Dahl et al. in the early 1960s [25, 26], is a widely used genetic model of human hypertension. The consomic rat strains, used as the normotensive control for the Dahl SS rat, are generated by substituting a chromosome or a part of a chromosome from a normal rat strain for the corresponding genomic region of the SS rat [24, 27, 28]. Previous research has shown that substitution of chromosome 13 or 18 can attenuate hypertension [31, 24]. Our study will focus on the hypertension-related genes listed in Ref. [24] by analysis of biological pathways.

¹ KEGG (Kyoto Encyclopedia of Genes and Genomes) website: <http://www.kegg.jp/> or <http://www.genome.jp/kegg/>.

Table 1: Examples of the correspondence between genes and pathways.

Gene	Gene ID	Pathway(s)
<i>Timp1</i>	116 510	rno04066
<i>Casp6</i>	83 584	rno04210
<i>Ank3</i>	361 833	rno05205
<i>Aqp1</i>	25 240	rno04964, rno04976
<i>Kcnj1</i>	24 521	rno04960, rno04971
<i>Ctsd</i>	171 293	rno04142, rno05152
<i>Hist1h2ai</i>	502 129	rno05034, rno05322
<i>Sdc1</i>	25 216	rno04512, rno04514, rno05144, rno05205
<i>Col4a1</i>	290 905	rno04151, rno04510, rno04512, rno04974, rno05146, rno05200, rno05222
<i>Fzd2</i>	64 512	rno04310, rno04390, rno04916, rno05166, rno05200, rno05205, rno05217

Let us consider an undirected network $G_H = (V_H, E_H)$, where $V_H = \{v_i\}$ ($i = 1, 2, \dots, N$) denotes the set of N nodes, and $E_H = \{v_i, v_j\}$ the set of edges or connections between nodes. We will use the following notation: $A_{ij} = 1$ indicates that there is an edge between nodes v_i and v_j ; and $A_{ij} = 0$ otherwise. Our pathway-based gene network model is constructed in two steps.

Step 1. Selection of genes according to pathways

We first extract nodes from the 335 different hypertension-related genes² given in Ref. [24] based on the KEGG PATHWAY Database. From this Database, we can easily obtain the information of whether a gene is involved in single or multiple pathways. Here, we only consider the pathways including hypertension-related genes, and the genes not involved in any pathway are excluded from our study. Thus, the $N = 90$ hypertension-related genes, which are involved in 157 pathways in the KEGG Database [32], are extracted to serve as nodes of the network model, where each node represents an individual gene shown as gene symbol or CloneID. Examples of ten genes involved in one or more pathways are shown in Table 1; here, a pathway is represented by the entry name of the Database, called the KEGG object identifier consisting of a database-dependent prefix and a five-digit number (such as rno04066, where the prefix “rno” designates the species to be rat).

Step 2. Establishment of connections

We now consider the correlations between any two genes according to the information of pathways. If two genes i and j are involved in the same pathway(s), then a connection is made between such two genes (nodes):

$$A_{ij} = \begin{cases} 1 & \text{if genes } i \text{ and } j \text{ are in the same pathway(s);} \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

Here $i, j = 1, 2, \dots, 90$ and $i \neq j$ (note that $A_{ii} = 0$ because we do not consider self-connections of nodes, i.e., self-interactions of genes). Consequently, there are biological regulatory relationships between genes i and j when $A_{ij} = 1$. In such a way, we have constructed the pathway-based network of hypertension-related genes, which contains 90 nodes (genes) and 482 edges (connections), as shown in Fig. 1.

3. Statistical and topological characteristics of the gene network

In this section, we analyze the pathway-based gene network of hypertension by calculating the following indices: degree (and average degree), degree distribution, average path length, clustering coefficient, assortativity coefficient, and four centrality indices (degree centrality, betweenness centrality, closeness centrality and integrated centrality), which can provide us with statistical and topological characteristics of the gene network.

3.1. Degree and degree distribution

The degree k_i of a node i is the number of edges connecting to the node. The average of k_i over all nodes is called the average degree of the network, and is denoted as $\langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i$. The degree distribution function $P(k)$ gives the probability that a randomly selected node has exactly k edges [3, 4]. The degree distribution is one of the most basic quantitative properties of a network.

² All these genes are given in Table S2 (Excel file) of Supplemental Figures and Tables of Ref. [24], which can be readily accessed from the website: <http://physiolgenomics.physiology.org/content/34/1/54/suppl/DC1>.

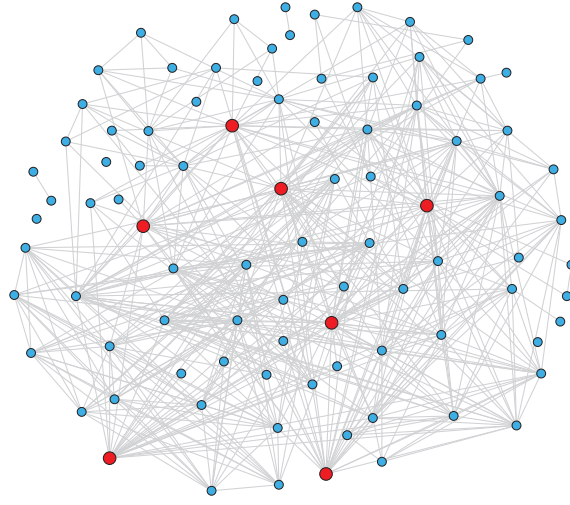


Figure 1: Illustration of the pathway-based gene network of hypertension for the SS rat with all 90 nodes and 482 edges. The large red nodes represent the seven hub genes. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

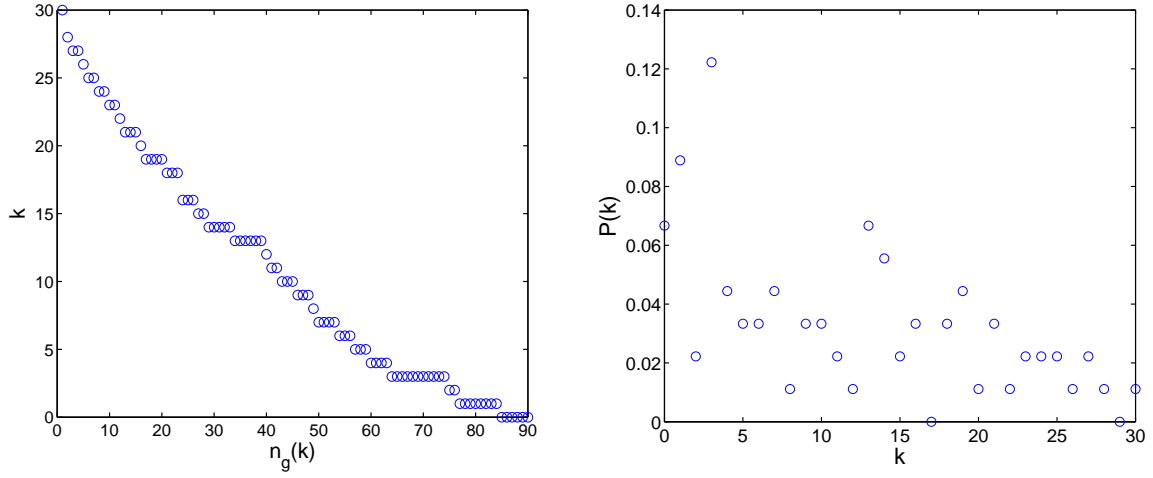


Figure 2: Degree and degree distribution of the pathway-based gene network: (left) The values of the degree of 90 nodes are shown in descending order; (right) The degree distribution is not scale-free.

For the pathway-based gene network of hypertension, the degree k_i of a node i is just the number of other genes which are involved in the same pathway(s) as gene i . Fig. 2 plots the degree and degree distribution of the gene network. Here, the values of the degree of 90 nodes are ranked in a descending order (with $n_g(k)$ denoting the rank of genes in degree values), and it is found that the gene *Jun* of $n_g = 1$ has the highest degree 30. The average degree $\langle k \rangle$ of this network is about 10.71. The illustration of degree distribution shows that the probability that a gene can link with k other genes does not decay as a power-law, suggesting that the gene network does not have a scale-free topology.

3.2. Average path length and clustering coefficient

In a network, a path from node i to node j is a sequence of adjacent nodes starting at i and ending at j . The path with the smallest number of edges between the two selected nodes is called the shortest path. The distance d_{ij} between two nodes i and j is defined as the number of edges along the shortest path connecting them. The diameter D is the maximum distance between any pair of nodes in the network, i.e., $D = \max\{d_{ij}\}$. The average path length L is defined as the mean distance between two nodes, averaged over all pairs of nodes, i.e.

$$L = \frac{1}{N(N-1)} \sum_{i \neq j} d_{ij}, \quad (2)$$

which offers a measure of the overall navigability of a network [33].

A clustering coefficient can be used to describe the cohesiveness of the neighborhood of a node [1, 4]. In a network, the clustering coefficient c_i is defined as the ratio between the number e_i of edges that actually link the k_i neighbors of node i to each other and the total possible number of edges among them, i.e.

$$c_i = \frac{2e_i}{k_i(k_i - 1)} \quad (k_i \geq 2). \quad (3)$$

The clustering coefficient C of the whole network is the average of c_i over all i , $C = \frac{1}{N} \sum_{i=1}^N c_i$, which characterizes the overall tendency of nodes to form clusters, clearly, $C \leq 1$.

The pathway-based gene network has a very short average path length: L is about 2.331, and $\log(\log N) < L < \log N$. The diameter D is only 5 (i.e., the distance between genes *Gcgr* and *Gnpat*), which implies at most five hops separate any two genes in the 482 connections of the network. The clustering coefficient is calculated to be $C = 0.6403$, which is relatively high. Therefore, the pathway-based gene network has the small-world property (short L and high C).

3.3. Assortativity

The concept of assortativity is introduced to describe degree correlations between neighboring nodes in a network [34]. A network is assortative if high (low) degree nodes tend to be connected to other high (low) degree nodes; otherwise, it is disassortative if high (low) degree nodes tend to be connected to other low (high) degree nodes. The assortativity can be described by the correlation between the degrees of neighboring nodes in terms of the mean Pearson correlation coefficient. Let x_i and y_i be the degrees of the end nodes of the i th edge, with $i = 1, 2, \dots, E$ (E is the number of edges in the network), then the assortativity coefficient of the network is given by [34]:

$$r = \frac{E^{-1} \sum_i x_i y_i - \left[E^{-1} \sum_i \frac{1}{2} (x_i + y_i) \right]^2}{E^{-1} \sum_i \frac{1}{2} (x_i^2 + y_i^2) - \left[E^{-1} \sum_i \frac{1}{2} (x_i + y_i) \right]^2}. \quad (4)$$

The network is assortative if $r > 0$, and disassortative if $r < 0$. The assortativity coefficient of the pathway-based gene network is calculated to be $r = 0.2178$, exhibiting an assortative behavior. This is different from most biological networks which show negative r .

3.4. Centrality

3.4.1. Definitions of four centrality indices

The degree centrality C_d , betweenness centrality C_b and closeness centrality C_c are three centrality indices commonly used in finding out the centralization nodes of the network [35, 36]. Recently, we also introduced an integrated centrality C_{integr} to fully reflect the contribution of three centrality indices $\{C_d, C_b, C_c\}$ [23].

The degree centrality of a given node i is the proportion of other nodes that are adjacent to node i [36], i.e.

$$C_d(i) = \frac{k_i}{N-1}, \quad (5)$$

Table 2: Top 25 values of degree centrality C_d , betweenness centrality C_b , closeness centrality C_c and integrated centrality C_{intgr} in the pathway-based gene network of hypertension. In this paper, nine genes are expressed as abbreviations (cf. Fig. 6 and its caption).

Gene	C_d	Gene	C_b	Gene	C_c	Gene	C_{intgr}
<i>Jun</i>	0.3371	<i>Sdhb</i>	0.09327	<i>Jun</i>	0.5087	<i>Jun</i>	0.9280
<i>Cdk4</i>	0.3146	<i>Jun</i>	0.07311	<i>Rps6kb1</i>	0.4880	<i>Rps6kb1</i>	0.8359
<i>RT1-Da</i>	0.3034	<i>RT1-Da</i>	0.06903	<i>Cycs</i>	0.4653	<i>Cycs</i>	0.8156
<i>Pdgfra</i>	0.3034	<i>Cycs</i>	0.06057	<i>Creb3l2</i>	0.4653	<i>Creb3l2</i>	0.8027
<i>Rps6kb1</i>	0.2921	<i>Rps6kb1</i>	0.05791	<i>Cdk4</i>	0.4446	<i>Cdk4</i>	0.7547
<i>Creb3l2</i>	0.2809	<i>Cd36</i>	0.05748	<i>Actg1</i>	0.4446	<i>Actg1</i>	0.7281
<i>Actg1</i>	0.2809	<i>Pdha1</i>	0.04875	<i>RT1-Da</i>	0.4413	<i>RT1-Da</i>	0.7128
<i>Csf1r</i>	0.2697	<i>Creb3l2</i>	0.04067	<i>Shc1</i>	0.4413	<i>Shc1</i>	0.6662
<i>Fn1</i>	0.2697	<i>Fcgr1</i>	0.03694	<i>Pdgfra</i>	0.4413	<i>Pdgfra</i>	0.6598
<i>Col4a1</i>	0.2584	<i>Ctsl</i>	0.03632	<i>Csf1r</i>	0.4381	<i>Csf1r</i>	0.6462
<i>Col4a2</i>	0.2584	<i>Cdk4</i>	0.03088	<i>Fn1</i>	0.4381	<i>Fn1</i>	0.6390
<i>Fzd2</i>	0.2472	<i>Shc1</i>	0.02745	<i>Cd36</i>	0.4318	<i>Cd36</i>	0.6293
<i>Cycs</i>	0.2360	<i>Csf1r</i>	0.02586	<i>Fzd2</i>	0.4287	<i>Fzd2</i>	0.6206
<i>Shc1</i>	0.2360	<i>Acs1l</i>	0.02133	<i>Col4a1</i>	0.4227	<i>Col4a1</i>	0.6041
<i>Tgfbr1</i>	0.2360	<i>Acs1l</i>	0.02133	<i>Col4a2</i>	0.4227	<i>Col4a2</i>	0.5956
<i>Fcgr1</i>	0.2247	<i>Pdgfra</i>	0.01974	<i>Tgfbr1</i>	0.4227	<i>Tgfbr1</i>	0.5841
<i>Col2a1</i>	0.2135	<i>Actg1</i>	0.01955	<i>Fcgr1</i>	0.4197	<i>Sdhb</i>	0.5678
<i>Col5a1</i>	0.2135	<i>Ghr</i>	0.01847	<i>Sdhb</i>	0.4111	<i>Fcgr1</i>	0.5678
<i>Col5a2</i>	0.2135	<i>Nrp1</i>	0.01839	<i>Mmp2</i>	0.4083	<i>Mmp2</i>	0.5640
<i>Col6a1</i>	0.2135	<i>Ctss</i>	0.01467	<i>Ctsl</i>	0.4056	<i>Ctsl</i>	0.5388
<i>Sdhb</i>	0.2022	<i>Cdc2a</i>	0.01256	<i>Col2a1</i>	0.4056	<i>Col2a1</i>	0.4980
<i>Mmp2</i>	0.2022	<i>Fzd2</i>	0.01080	<i>Col5a1</i>	0.4056	<i>Col5a1</i>	0.4980
<i>Ctsl</i>	0.2022	<i>Polr3e</i>	0.01018	<i>Col5a2</i>	0.4056	<i>Col5a2</i>	0.4980
<i>Cd36</i>	0.1798	<i>Col4a1</i>	0.00987	<i>Col6a1</i>	0.4056	<i>Col6a1</i>	0.4980
<i>Pdha1</i>	0.1798	<i>Col4a2</i>	0.00987	<i>Sdc1</i>	0.3949	<i>Sdc1</i>	0.4796

here $N - 1$ is the maximum possible degree of the network.

The betweenness centrality of a node i is defined as the proportion of all shortest paths (geodesics) between pairs of other nodes that include this node i [36]:

$$C_b(i) = \sum_{j(<k)} \sum_k \frac{g_{jk}(i)}{g_{jk}}, \quad (6)$$

where g_{jk} is the number of shortest paths between nodes j and k , and $g_{jk}(i)$ the number of shortest paths containing node i between nodes j and k .

The closeness centrality of a node i is the number of other nodes divided by the sum of the distances between node i and all others [35, 36]:

$$C_c(i) = (L_i)^{-1} = \frac{N - 1}{\sum_{j=1}^N d_{ij}}, \quad (7)$$

here L_i is the average distance between node i and all other nodes.

To comprehensively and quantitatively reflect the contribution of the above three centrality indices $\{C_d, C_b, C_c\}$, we can also introduce the integrated centrality C_{intgr} of node i , defined as follows [23]:

$$C_{\text{intgr}}(i) = \frac{1}{3} \left[\frac{C_d(i)}{C_{d,\max}} + \frac{C_b(i)}{C_{b,\max}} + \frac{C_c(i)}{C_{c,\max}} \right], \quad (8)$$

where $C_{d,\max}$, $C_{b,\max}$ and $C_{c,\max}$ are the maximums of $\{C_d\}$, $\{C_b\}$ and $\{C_c\}$, respectively. Obviously, $C_{\text{intgr}}(i)$ has a value between 0 and 1.

3.4.2. Centrality analysis and hub genes

For the pathway-based gene network of hypertension, the four centrality indices are calculated based on the above definitions, and the top 25 values of each centrality index are listed in Table 2. In the following, we will determine hub genes in the network through centrality analysis.

Fig. 3 describes the correspondence among degree centrality C_d , betweenness centrality C_b and closeness centrality C_c of nodes in the pathway-based gene network. We can see in Fig. 3 that C_c of most of the nodes is distributed

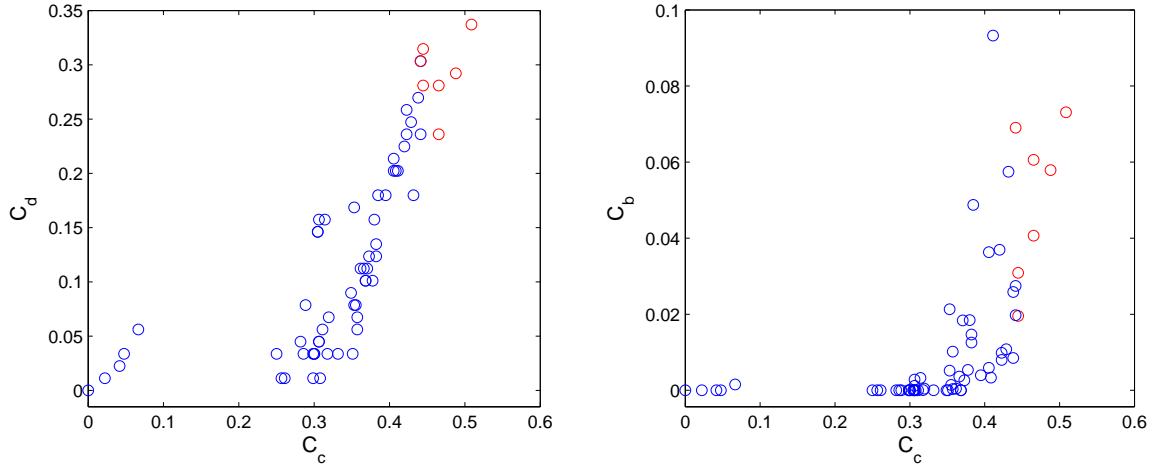


Figure 3: Correspondence among degree centrality C_d , betweenness centrality C_b and closeness centrality C_c of nodes in the pathway-based gene network: (left) C_d versus C_c ; (right) C_b versus C_c . The red circles represent the seven hub genes. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

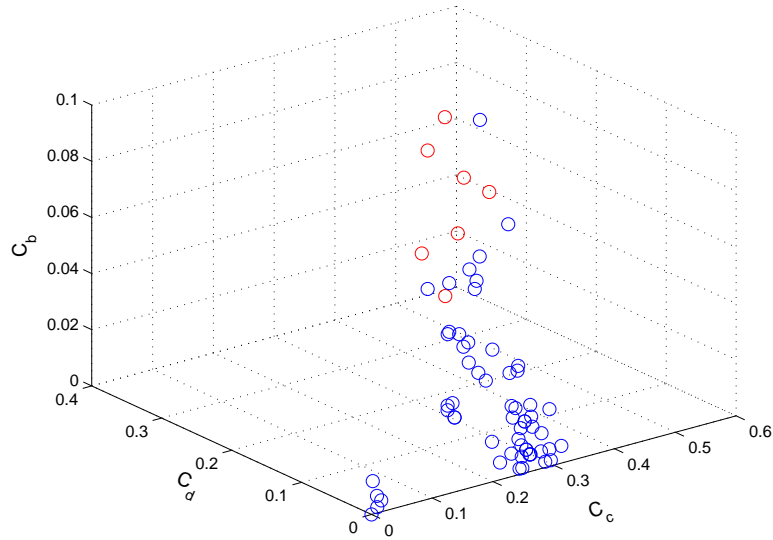


Figure 4: Three-dimensional diagram of degree centrality C_d , betweenness centrality C_b and closeness centrality C_c of nodes in the pathway-based gene network. The red circles represent the seven hub genes. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

in the range of (0.25, 0.51). In this range of C_c , we also observe that in Fig. 3 (left) the distribution of C_d is roughly uniform in the range of (0.01, 0.34); while in Fig. 3 (right) C_b is non-uniformly distributed in the range of [0, 0.094), here most of the nodes have very small C_b , and only a few nodes have large C_b which also have relatively large C_c . Fig. 4 combines the three centrality indices $\{C_d, C_b, C_c\}$ in the three-dimensional space. It can be seen from Fig. 4 that a small number of nodes have high values of three centrality indices, which can be viewed as hubs of the network. Generally speaking, each of these three centrality indices has its own focus on an influence of a node on other nodes in the network, thus one can identify hubs according to that focus. However, in order to fully reflect the contribution of all these three centrality indices, we will simply determine key hub genes using integrated centrality C_{integr} , as in Ref. [23].

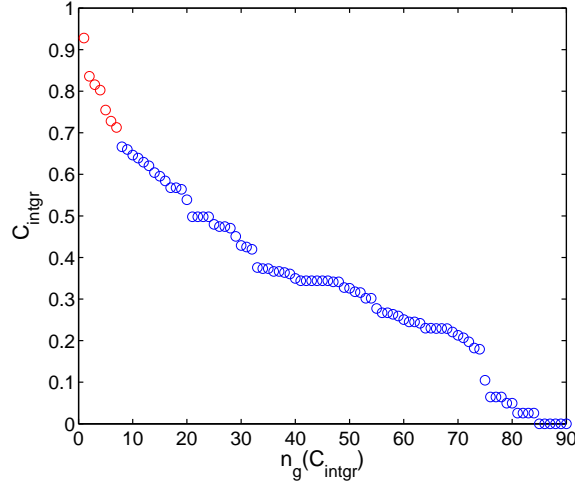


Figure 5: The values of integrated centrality C_{intgr} of 90 nodes in descending order for the pathway-based gene network. The red circles represent the seven hub genes. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Table 3: Biological descriptions of seven hub genes in the pathway-based gene network.

Gene	Official full name	Gene ID	Description
<i>Jun</i>	jun proto-oncogene	24 516	Encodes a protein that exhibits double-stranded DNA binding
<i>Rps6kb1</i>	ribosomal protein S6 kinase, polypeptide 1	83 840	Encodes a protein that exhibits ATP binding; peptide binding
<i>Cyts</i>	cytochrome c, somatic	25 309	Encodes a protein that exhibits electron carrier activity
<i>Creb3l2</i>	cAMP responsive element binding protein 3-like 2	362 339	Encodes a protein that exhibits cAMP response element binding
<i>Cdk4</i>	cyclin-dependent kinase 4	94 201	Related to circadian rhythm; organ regeneration
<i>Actg1</i>	actin, gamma 1	287 876	Involved in response to calcium ion
<i>RT1-Da</i>	RT1 class II, locus Da	294 269	Involved in antigen processing

As calculated and shown above, the pathway-based gene network is not scale-free. Fig. 5 illustrates the integrated centrality C_{intgr} of 90 nodes in a descending order (with $n_g(C_{\text{intgr}})$ denoting the rank of genes in C_{intgr}) for the network, which shows that the top seven nodes (red) have high values of C_{intgr} . From Table 2 and Fig. 5, seven nodes are determined as important hub genes: *Jun*, *Rps6kb1*, *Cyts*, *Creb3l2*, *Cdk4*, *Actg1* and *RT1-Da*. These seven genes have the integrated centrality of $C_{\text{intgr}} > 0.71$, which is about 77% of its maximum (0.9280). However, it should be mentioned that there is no strict significance threshold for C_{intgr} , and one can also lower the threshold to enable more genes to be included in hub genes.

In Table 3 we list the official full names, gene IDs and biological descriptions of seven hub genes. Here, we also provide a brief biological description of the first hub gene. The gene *Jun* is ranked first in C_{intgr} because it has the highest C_d and C_c , as well as the second highest C_b . *Jun* encodes a protein that exhibits double-stranded DNA binding, involved in aging, angiogenesis, endothelin and Rho/Rac/Cdc42 mediated signaling pathways. It is associated with kidney neoplasms and spinal cord injuries. Recent studies have confirmed that *Jun* is closely related to low potassium reaction and cell renal cell carcinoma. A low potassium diet might induce hypertension, which is always accompanied by hypokalemia. The incidence of renal cell carcinoma coupled with hypertension is up to 14%–40% [37].

Among the seven hub genes, we note that the gene *Cdk4* is the only common hub gene in both the pathway-based gene network here and the gene co-expression network of our previous study [23]. Besides *Cdk4*, we also see in Table 2 that the three genes *Shc1*, *Fzd2* and *Col4a1*, which have relatively high integrated centrality $C_{\text{intgr}} > 0.60$, are the hub genes identified in the gene co-expression network of Ref. [23]. These four genes have been confirmed by biological and medical research to play important roles in hypertension. Moreover, we also observe that although the gene *Sdhb* is ranked only joint 17th in C_{intgr} ($= 0.5678$), it has the highest betweenness centrality C_b . Since C_b is based on the shortest paths and reflects the ability of a node to influence other related nodes in the network, *Sdhb* should also be a key gene in hypertension. If we lower the threshold of integrated centrality to $C_{\text{intgr}} > 0.50$, then these four genes (*Shc1*, *Fzd2*, *Col4a1* and *Sdhb*) can also be included in hub genes. In this paper, we do not take this

Table 4: Dissimilarity scores $d_s(i, j)$ of ten nodes (genes) in the pathway-based gene network. The specific modules, which are identified after the completion of the modular decomposition, are indicated in parentheses in the first column.

	<i>Jun</i>	<i>RT1-Da</i>	<i>Col4a1</i>	<i>Fzd2</i>	<i>Sdhb</i>	<i>Gda</i>	<i>Sec13l1</i>	<i>Sumo1</i>	<i>Aqp1</i>	<i>Kcnj1</i>
<i>Jun</i> (I)	0.0000	0.4310	0.2931	0.2414	0.7241	0.6897	0.6034	0.5517	0.5345	0.5172
<i>RT1-Da</i> (I)	0.4310	0.0000	0.5172	0.3621	0.7069	0.7069	0.5517	0.5000	0.4828	0.4655
<i>Col4a1</i> (I)	0.2931	0.5172	0.0000	0.2931	0.6034	0.6379	0.4828	0.4310	0.4138	0.3966
<i>Fzd2</i> (I)	0.2414	0.3621	0.2931	0.0000	0.5862	0.6207	0.4655	0.4138	0.3966	0.3793
<i>Sdhb</i> (II)	0.7241	0.7069	0.6034	0.5862	0.0000	0.1034	0.3966	0.3448	0.3276	0.3103
<i>Gda</i> (II)	0.6897	0.7069	0.6379	0.6207	0.1034	0.0000	0.3276	0.2759	0.2586	0.2414
<i>Sec13l1</i> (III)	0.6034	0.5517	0.4828	0.4655	0.3966	0.3276	0.0000	0.0517	0.1034	0.0862
<i>Sumo1</i> (III)	0.5517	0.5000	0.4310	0.4138	0.3448	0.2759	0.0517	0.0000	0.0517	0.0345
<i>Aqp1</i> (IV)	0.5345	0.4828	0.4138	0.3966	0.3276	0.2586	0.1034	0.0517	0.0000	0.0172
<i>Kcnj1</i> (V)	0.5172	0.4655	0.3966	0.3793	0.3103	0.2414	0.0862	0.0345	0.0172	0.0000

lower threshold of $C_{\text{integr}} > 0.50$ because, in view of the relatively large range of variation of C_b , we do not want the genes of small C_b to be included in hub genes. In a wider view, however, these four genes (*Shc1*, *Fzd2*, *Col4a1* and *Sdhb*), together with the above seven hub genes, are worthy of further study in the future.

4. Modular structure of the gene network

Many networks are found to divide naturally into modules or communities, i.e., groups of nodes within which the connections are relatively dense but between which they are sparser [38, 39]. In this section, we explore the modular structure of the pathway-based gene network of hypertension.

4.1. Structural equivalence of nodes

Two nodes are structural equivalent if they have identical connections with all other nodes. We can use a dissimilarity index d_s to measure the equivalence of two nodes i and j as follows [40]:

$$d_s(i, j) = \frac{|V(i) + V(j)|}{k_1 + k_2}. \quad (9)$$

Here $i, j = 1, 2, \dots, 90$, $V(i)$ are all neighbors of node i , $|\cdot|$ stands for set cardinality, k_1 and k_2 stand for the largest and the second largest degree in the network, respectively. Obviously, $d_s(i, j)$ has a value between 0 (completely similar) and 1 (completely different). In Table 4 we list the dissimilarity scores of ten genes distributed in different modules (cf. Fig. 7).

The dissimilarity scores allow us to cluster nodes in accordance with the structural equivalence into the corresponding positions by the hierarchical clustering technique. First, the nodes that are most similar are grouped into a cluster. Then, the next pair of nodes or clusters that are most similar are grouped, and this process continues until all nodes have been joined. The dendrogram in Fig. 6 is obtained with Pajek software [40, 41], which visualizes the above clustering process.

4.2. Modular decomposition of the network

Based on the above dendrogram, we can obtain the modular structure of the gene network, which is shown in Fig. 7. The network consists of five modules:

- (I) the largest module with 58 nodes (red and green);
- (II) the second largest module with 16 nodes (blue);
- (III) a small module with the highest clustering coefficient including 6 nodes (brown);
- (IV) two pairs of adjacent nodes (each joined by a single connection) (purple);
- (V) six isolated nodes (yellow).

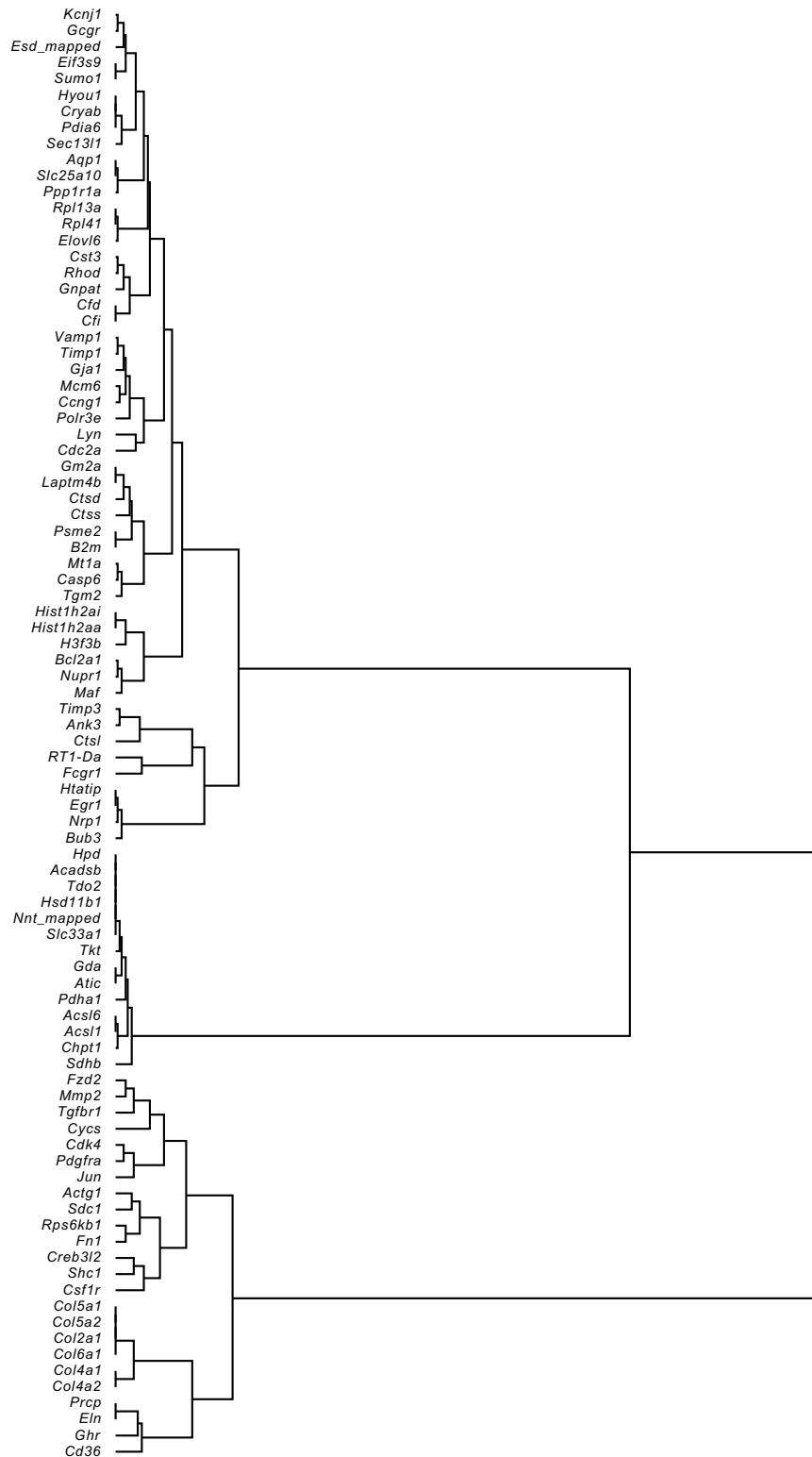


Figure 6: Dendrogram of similarities. It shows the hierarchical clustering of the pathway-based gene network. In this paper, the following nine genes (as recorded in Ref. [24]), *Rhod*-predicted, *Polr3e*-predicted, *Hist1h2ai*-predicted // *Hist1h4a*-predicted, *Sdhb*-predicted, *Actg1* // *LOC295810*, *Col6a1*-predicted, *Col4a2*-predicted, *Prpc*-predicted, and *Cd36* // *RGD1562323*-predicted, are abbreviated as *Rhod*, *Polr3e*, *Hist1h2ai*, *Sdhb*, *Actg1*, *Col6a1*, *Col4a2*, *Prpc*, and *Cd36*, respectively.

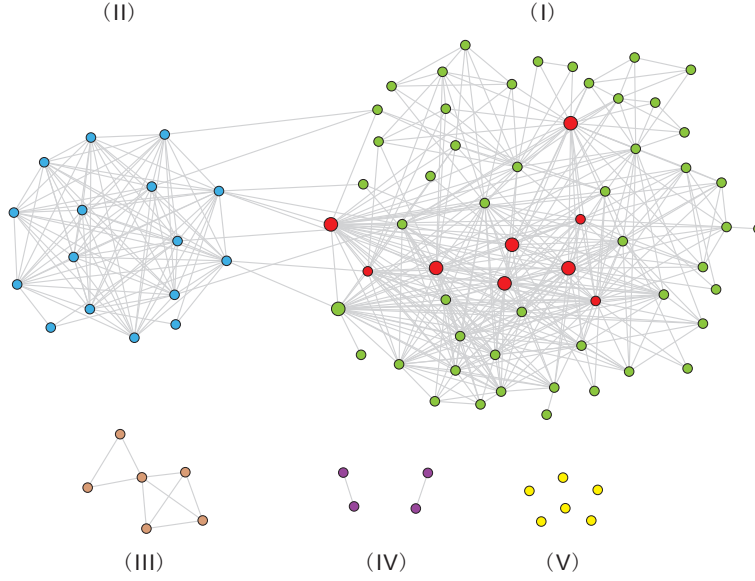


Figure 7: Modular structure of the pathway-based gene network. Modules I–V contain 58, 16, 6, 4 and 6 nodes, respectively. The nine nodes marked red in module I correspond to the top nine genes involved in the largest number of pathways. The seven large nodes (six red and one green) indicate the hub genes. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Table 5: Statistical characteristics of five modules in the pathway-based gene network.

Module	N_m	N_m/N	$\langle k \rangle_m$	C_m
I	58	64.44%	12.72	0.6545
II	16	17.78%	12.75	0.8917
III	6	6.67%	3	0.9000
IV	4	4.44%	1	0
V	6	6.67%	0	0

We can observe that modules I and II constitute the largest connected part of the network, containing 74 nodes and 471 edges; and the nodes are highly connected within a module but much less connected between modules.

In Table 5 we list several statistical characteristics of every module, including the number N_m of nodes, node proportion N_m/N , average degree $\langle k \rangle_m$ and clustering coefficient C_m ($m = \text{I, II, III, IV, V}$). Here C_m is the average of c_i over all nodes in module m . From Table 5 we can see that the average degree of each of the two modules I and II ($\langle k \rangle_{\text{I}} = 12.72$, $\langle k \rangle_{\text{II}} = 12.75$) is greater than the average degree $\langle k \rangle = 10.71$ of the whole network. Except modules IV and V of $C_{\text{IV}} = C_{\text{V}} = 0$, the clustering coefficient of each of other three modules (I–III) exceeds that of the whole network ($C = 0.6403$), showing that there are more connections within each of these three modules, which justifies the modular decomposition of the whole network.

4.3. Characteristics of modular structure

The connections between nodes in the gene network are created based on whether genes are involved in the same pathway(s). Having examined the number of pathways in which each gene is involved, we find that there are nine genes involved in more than 10 pathways: *Jun* (30), *RT1-Da* (21), *Cdk4* (17), *Cycs* (17), *Creb3l2* (16), *Actg1* (15), *Pdgfra* (14), *Tgfbr1* (14) and *Shc1* (13), here the number of pathways involved is shown in parentheses. These nine genes (red nodes in Fig. 7) are a basis of dense connections within the largest module I. We also note that among these nine genes, there are six hub genes (large nodes marked red); another hub gene *Rps6kb1* (large node marked green) is involved in nine pathways.

We can examine the robustness of the network based on modular structure [39, 42]. The nodes within a module (except IV and V in this paper) are relatively robust against mutation because there are multiple paths between any two nodes and thus the network will not be easily broken when mutation occurs. Nevertheless, the parts of fewer

connections between modules should be the weaknesses of the network system. Fig. 7 visualizes the weak links of the network. The removal of these weak connections and relevant nodes would result in the breaking of the network. Thus, the modular structure analysis can facilitate the exploration of the relationship between weak connections of the network and drug targets of hypertension.

The genes in the largest connected part (i.e., core modules I and II) are involved in multiple pathways, and contribute to a variety of biological functions. It is difficult to assign each of these genes to a single biological function because of pleiotropy, namely, one gene might influence many different biological processes in organisms [43]. However, the genes in the small non-core module III with six nodes are involved in only two specific pathways, rno03013 (RNA transport) and rno04141 (protein processing in endoplasmic reticulum (ER)), which indicates that the non-core module III corresponds to clearly identified pathways and functions relatively independently.

In the pathway rno03013, the different RNA species produced in the nucleus are exported through the nuclear pore complexes (NPCs) to the cytoplasm via mobile export receptors, which is fundamental for gene expression. In the pathway rno04141, newly synthesized peptides enter the ER via the sec61 pore and are glycosylated; correctly folded proteins are packaged into transport vesicles and misfolded proteins are retained within the ER lumen (cf. KEGG PATHWAY Database).

5. Summary and concluding remarks

Hypertension is a cardiovascular disease associated with long-term interaction between genetic and environmental factors. In this study, we use pathways data to obtain backwards the relationships between the hypertension-related genes, try to extract the complex interactions between genes through calculating statistical characteristics and analyzing modular structure of the network.

The pathway-based gene network has the following characteristics: (i) The network does not obey a power-law degree distribution and thus is not of a scale-free property. The seven hub genes that are identified by integrated centrality $C_{\text{integr}} > 0.71$ are: *Jun*, *Rps6kb1*, *Cycs*, *Creb3l2*, *Cdk4*, *Actg1* and *RT1-Da*; they are key (feature) genes involved in the formation of hypertension. (ii) The network shows the small-world property (i.e., a small L and a large C), which reveals the direct influence of these hub genes on hypertension from another perspective. (iii) The network has a modular structure. The weak connections of the network can be visualized by its modular structure, which can help to screen out key hypertension-related genes or pathways.

In this paper, we construct the network model of hypertension-related genes based on biological pathways. Among the seven hub genes identified in this network, only *Cdk4* is also a hub gene in the gene co-expression network of our previous study [23]. Besides *Cdk4*, the three genes *Shc1*, *Fzd2* and *Col4a1* with $C_{\text{integr}} > 0.60$ in the pathway-based gene network are identified as the hub genes in the gene co-expression network of Ref. [23]. Although we can see that more nodes will become hub genes and thus there will be more hub genes overlapped in the both networks if we lower the threshold of C_{integr} , the hub genes in the two networks are impossible to be completely overlapped because the two networks are constructed from the different perspectives, i.e., based on the different characters of the genes. The results from this paper and the theoretical analysis in Ref. [23] would complement each other. The seven hub genes in the pathway-based gene network, together with the above-mentioned *Shc1*, *Fzd2* and *Col4a1*, as well as *Sdhb* of the highest C_b , can be regarded as candidate genes or drug targets for further biological and medical research on their functions; in particular, the common hub gene *Cdk4* of the both networks would be worth more attention.

Moreover, the network may also be analyzed based on other functional correlation methods, such as GO analysis [21], to get more molecular mechanisms about hypertension. In the next study, we will develop weighted network models and explore the mutual regulatory relationships between genes of complex diseases using dynamical analysis. These studies will provide another perspective on expounding the differentially expressed genes and finding new drug targets for other serious diseases. Finally, we expect that the complex network approach can provide clues for exploring the pathogenesis of critical illness from molecular perspective.

Acknowledgments

This work was supported in part by the National Natural Science Foundation of China (NSFC) (Grant Nos. 11365023 and 61263043), the Projects of the Science and Technology Research and Development Program of Baoji

City (Grant Nos. 15RKX-1-5-14 and 15RKX-1-5-6), the Key Project of Baoji University of Arts and Sciences (Grant No. ZK14035), and the Joint Fund of Department of Science and Technology of Guizhou Province, Bureau of Science and Technology of Qiandongnan Prefecture, and Kaili University (Grant No. LH-2014-7231). We are grateful to the authors of Ref. [24] for providing the gene information of the SS rat, which is the basis of construction of our network model. The authors would like to thank Professor Huai Cao for his helpful discussions and suggestions.

References

- [1] D.J. Watts, S.H. Strogatz, Collective dynamics of 'small-world' networks, *Nature* 393 (6684) (1998) 440–442.
- [2] A.-L. Barabási, R. Albert, Emergence of scaling in random networks, *Science* 286 (5439) (1999) 509–512.
- [3] S.H. Strogatz, Exploring complex networks, *Nature* 410 (6825) (2001) 268–276.
- [4] M.E.J. Newman, The structure and function of complex networks, *SIAM Rev.* 45 (2) (2003) 167–256.
- [5] W.-X. Wang, C.-Y. Yin, G. Yan, B.-H. Wang, Integrating local static and dynamic information for routing traffic, *Phys. Rev. E* 74 (1) (2006) 016101.
- [6] H. Jeong, B. Tombor, R. Albert, Z.N. Oltvai, A.-L. Barabási, The large-scale organization of metabolic networks, *Nature* 407 (6804) (2000) 651–654.
- [7] A.-L. Barabási, Z.N. Oltvai, Network biology: Understanding the cell's functional organization, *Nature Rev. Genet.* 5 (2) (2004) 101–113.
- [8] F. Karlsson, M. Hörnquist, Order or chaos in Boolean gene networks depends on the mean fraction of canalizing functions, *Physica A* 384 (2) (2007) 747–757.
- [9] M. Tsuchiya, K. Selvarajoo, V. Piras, M. Tomita, A. Giuliani, Local and global responses in complex gene regulation networks, *Physica A* 388 (8) (2009) 1738–1746.
- [10] L. Diambra, Coarse-grain reconstruction of genetic networks from expression levels, *Physica A* 390 (11) (2011) 2198–2207.
- [11] O.A. Carretero, S. Oparil, Essential hypertension: Part I: Definition and etiology, *Circulation* 101 (3) (2000) 329–335.
- [12] A.W. Cowley Jr., J.H. Nadeau, A. Baccarelli, K. Berecek, M. Fornage, G.H. Gibbons, D.G. Harrison, M. Liang, P.W. Nathanielsz, D.T. O'Connor, J. Orдовas, W. Peng, M.B. Soares, M. Szyf, H.E. Tolunay, K.C. Wood, K. Zhao, Z.S. Galis, Report of the National Heart, Lung, and Blood Institute Working Group on epigenetics and hypertension, *Hypertension* 59 (5) (2012) 899–905.
- [13] M.H. Weinberger, N.S. Fineberg, S.E. Fineberg, M. Weinberger, Salt sensitivity, pulse pressure, and death in normal and hypertensive humans, *Hypertension* 37 (2) (2001) 429–432.
- [14] M.H. Alderman, Salt, blood pressure and health: A cautionary tale, *Int. J. Epidemiol.* 31 (2) (2002) 311–315.
- [15] B. Rodriguez-Iturbe, N.D. Vaziri, Salt-sensitive hypertension—update on novel findings, *Nephrol. Dial. Transplant.* 22 (4) (2007) 992–995.
- [16] J.P. Rapp, Genetic analysis of inherited hypertension in the rat, *Physiol. Rev.* 80 (1) (2000) 135–172.
- [17] A.W. Cowley Jr., The genetic dissection of essential hypertension, *Nature Rev. Genet.* 7 (11) (2006) 829–840.
- [18] R. Cooper, J. Cutler, P. Desvigne-Nickens, S.P. Fortmann, L. Friedman, R. Havlik, G. Hogelin, J. Marler, P. McGovern, G. Morosco, L. Mosca, T. Pearson, J. Stamler, D. Stryer, T. Thom, Trends and disparities in coronary heart disease, stroke, and other cardiovascular diseases in the United States: Findings of the National Conference on Cardiovascular Disease Prevention, *Circulation* 102 (25) (2000) 3137–3147.
- [19] K. Wolf-Maier, R.S. Cooper, H. Kramer, J.R. Banegas, S. Giampaoli, M.R. Joffres, N. Poulter, P. Primatesta, B. Stegmayr, M. Thamm, Hypertension treatment and control in five European countries, Canada, and the United States, *Hypertension* 43 (1) (2004) 10–17.
- [20] H. Sanada, J.E. Jones, P.A. Jose, Genetics of salt-sensitive hypertension, *Curr. Hypertens. Rep.* 13 (1) (2011) 55–66.
- [21] F. Censi, A. Giuliani, P. Bartolini, G. Calcagnini, A multiscale graph theoretical approach to gene regulation networks: A case study in atrial fibrillation, *IEEE Trans. Biomed. Eng.* 58 (10) (2011) 2943–2946.
- [22] R. Demicheli, D. Coradini, Gene regulatory networks: A new conceptual framework to analyse breast cancer behaviour, *Ann. Oncol.* 22 (6) (2011) 1259–1265.
- [23] H. Wang, C.-Y. Xu, J.-B. Hu, K.-F. Cao, A complex network analysis of hypertension-related genes, *Physica A* 394 (2014) 166–176.
- [24] M. Liang, N.H. Lee, H. Wang, A.S. Greene, A.E. Kwitek, M.L. Kaldunski, T.V. Luu, B.C. Frank, S. Bugenhagen, H.J. Jacob, A.W. Cowley Jr., Molecular networks in Dahl salt-sensitive hypertension based on transcriptome analysis of a panel of consomic rats, *Physiol. Genomics* 34 (1) (2008) 54–64.
- [25] L.K. Dahl, M. Heine, L. Tassinari, Effects of chronic excess salt ingestion: Evidence that genetic factors play an important role in susceptibility to experimental hypertension, *J. Exp. Med.* 115 (6) (1962) 1173–1190.
- [26] J.P. Rapp, Dahl salt-susceptible and salt-resistant rats: A review, *Hypertension* 4 (6) (1982) 753–763.
- [27] A.Y. Deng, In search of hypertension genes in Dahl salt-sensitive rats, *J. Hypertens.* 16 (12) (1998) 1707–1717.
- [28] A.W. Cowley Jr., R.J. Roman, H.J. Jacob, Application of chromosomal substitution techniques in gene-function discovery, *J. Physiol.* 554 (1) (2004) 46–55.
- [29] H. Kitano, Systems biology: A brief overview, *Science* 295 (5560) (2002) 1662–1664.
- [30] Y. Deville, D. Gilbert, J. van Helden, S.J. Wodak, An overview of data models for the analysis of biochemical pathways, *Brief. Bioinform.* 4 (3) (2003) 246–259.
- [31] B. Joe, N.E. Letwin, M.R. Garrett, S. Dhindaw, B. Frank, R. Sultana, K. Verratti, J.P. Rapp, N.H. Lee, Transcriptional profiling with a blood pressure QTL interval-specific oligonucleotide array, *Physiol. Genomics* 23 (3) (2005) 318–326.
- [32] H. Wang, Analyses of hypertension-related genes based on complex network theory (Doctor of Science Dissertation), Yunnan University, Kunming, 2013 (in Chinese, with English abstract).
- [33] J. Xu, Theory and Application of Graphs, Kluwer Academic Publishers, Dordrecht, Boston, London, 2003.
- [34] M.E.J. Newman, Assortative mixing in networks, *Phys. Rev. Lett.* 89 (20) (2002) 208701.
- [35] M.A. Beauchamp, An improved index of centrality, *Behav. Sci.* 10 (2) (1965) 161–163.
- [36] L.C. Freeman, Centrality in social networks: Conceptual clarification, *Social Networks* 1 (3) (1978–1979) 215–239.

- [37] E. Konik, E.G. Kurtz, F. Sam, D. Sawyer, Coronary artery spasm, hypertension, hypokalemia and licorice, *J. Clin. Case Rep.* 2 (8) (2012) 143.
- [38] M. Girvan, M.E.J. Newman, Community structure in social and biological networks, *Proc. Natl. Acad. Sci. USA* 99 (12) (2002) 7821–7826.
- [39] M.E.J. Newman, Modularity and community structure in networks, *Proc. Natl. Acad. Sci. USA* 103 (23) (2006) 8577–8582.
- [40] W. de Nooy, A. Mrvar, V. Batagelj, Exploratory Social Network Analysis with Pajek, in: *Structural Analysis in the Social Sciences*, vol. 34, Cambridge University Press, New York, 2005; Revised and expanded second edition, 2011.
- [41] V. Batagelj, A. Mrvar, Pajek: A program for large network analysis, *Connections* 21 (2) (1998) 47–57.
- [42] D.S. Callaway, M.E.J. Newman, S.H. Strogatz, D.J. Watts, Network robustness and fragility: Percolation on random graphs, *Phys. Rev. Lett.* 85 (25) (2000) 5468–5471.
- [43] W.G. Hill, X.-S. Zhang, On the pleiotropic structure of the genotype–phenotype map and the evolvability of complex organisms, *Genetics* 190 (3) (2012) 1131–1137.